

Original Article

# Multi-Objective Optimization Models for Performance, Energy, and Security in Hybrid Cloud Infrastructures

\*Amaka Udo<sup>1</sup>, Vo Thi Mai<sup>2</sup>, Amani Abdelrahman<sup>3</sup>

<sup>1,2,3</sup> Faculty of Data Science, Enugu State University of Science and Technology, Enugu, Nigeria.

## Abstract:

Hybrid clouds spanning on-premises; edge; and multiple public providers demand orchestration that balances conflicting goals: low latency and high throughput; energy and carbon reduction; and rigorous security. This work frames resource provisioning; workload placement; and autoscaling as a multi-objective optimization problem under uncertainty. We model application components (VMs/containers/functions) with performance SLOs; security postures; and energy carbon profiles tied to time-varying grid intensity and datacenter PUE. Objectives minimize end-to-end latency and SLO violations; total energy and carbon cost; and security risk (e.g.; exposure time; vulnerability risk scores; data-in-motion footprint); while respecting budget; data-sovereignty; and trust-zone constraints. Solution strategies combine exact formulations (MILP for small horizons) with scalable heuristics and metaheuristics (e.g.; NSGA-II/MOEA-D) to approximate the Pareto frontier; and a model-predictive layer that adapts to demand; spot price volatility; and carbon signals. Robust and risk-aware variants incorporate chance constraints and CVaR to hedge against workload and failure uncertainty. We integrate zero-trust placement rules; encryption overhead; and privacy controls (federated learning/differential privacy) as first-class constraints; and co-optimize network paths to limit cross-cloud data exfiltration risk. A learning-augmented scheduler uses surrogate models for queueing delays and power draw to accelerate search at runtime. The resulting policies enable operators to trade milliseconds of tail latency for double-digit energy/carbon savings or reduced attack surface; making decisions transparent via Pareto sets and what-if analysis. The framework generalizes across microservices; data pipelines; and AI inference; and can plug into Kubernetes-centric control planes.

## Keywords:

Hybrid Cloud, Multi-Objective Optimization, Pareto Frontier, NSGA-II, MOEA/D, Model-Predictive Control, Robust Optimization, Cvar Risk, Carbon-Aware Scheduling, Energy Efficiency, SLO-Aware Autoscaling, Secure Workload Placement, Zero-Trust Architecture, Data-Sovereignty Constraints, Differential Privacy, Federated Learning, Queueing-Theoretic Surrogates, Kubernetes Orchestration, Edge Cloud Continuum.

## Article History:

Received: 22.01.2022

Revised: 07.02.2022

Accepted: 16.02.2022

Published: 03.03.2022

## 1. Introduction

Hybrid cloud infrastructures spanning on-premises datacenters, edge nodes, and multiple public clouds have become the default substrate for modern digital services. This heterogeneity offers elastic capacity, proximity to users, and specialized accelerators, but it also introduces conflicting operational goals. Performance teams seek low end-to-end latency and high throughput under stochastic



demand; sustainability teams target energy and carbon reductions amid variable power usage effectiveness (PUE) and grid carbon intensity; security teams enforce zero-trust principles, data-sovereignty rules, and encryption with measurable overheads. Traditional single-objective controllers (e.g., minimize cost or latency alone) are brittle in this setting: they violate service-level objectives (SLOs) during bursts, overspend energy budgets, or expand the attack surface via risky cross-cloud data movement. The resulting decision landscape is inherently multi-objective and dynamic, with uncertainty arising from workload arrivals, resource contention, spot price volatility, failures, and evolving threats.

This paper frames workload placement, resource provisioning, and autoscaling as a multi-objective optimization problem that simultaneously considers performance, energy/carbon, and security risk. We model microservices, data pipelines, and AI inference components with queueing-based latency surrogates, power and carbon models tied to location and time, and security exposure metrics (e.g., data-in-motion footprint, vulnerability risk scores, and blast radius). We combine exact formulations for short horizons with scalable evolutionary search (NSGA-II/MOEA-D) and a learning-augmented, model-predictive controller that updates decisions as signals change. Robust and risk-aware variants encode chance constraints and CVaR to hedge against tail events. The key contributions are: (i) a unified objective and constraint model integrating zero-trust policies and sovereignty rules; (ii) fast surrogate models enabling online Pareto exploration; and (iii) actionable trade-off frontiers that let operators exchange milliseconds of tail latency for double-digit energy/carbon savings or reduced attack surface, all within Kubernetes-centric control planes.

## 2. Related Work

### 2.1. Performance Optimization Approaches in Cloud Systems

Classic performance engineering in clouds spans queueing-theoretic autoscalers, SLO-aware admission control, and topology-aware placement. Early work used M/M/1 and M/G/1 surrogates to convert latency targets into replica counts and CPU caps, while cost-based schedulers minimized instance-hours under throughput constraints. With microservices, researchers explored co-location/anti-affinity rules and interference-aware placement using contention models for CPU caches, memory bandwidth, and noisy neighbors. Recent systems add tail-latency control (p95/p99) via hedged requests, request reissuing, and adaptive timeouts, often coupled with service meshes to steer traffic along low-queue paths. Learning-augmented schedulers employ predictive signals (ARIMA/LSTM) for arrival rates and use MPC to pre-warm capacity or provision burstable instances, reducing cold-start penalties for serverless and GPU-backed inference.

### 2.2. Energy Efficiency and Green Computing Techniques

Energy work spans hardware-, platform-, and workload-level levers. At the hardware layer, DVFS and power capping trade frequency for watts; at the platform layer, consolidation and sleep states (C-states) pack workloads to power down hosts while watching for SLA regressions. Carbon-aware scheduling extends energy minimization by aligning jobs with diurnal grid-carbon intensity and location-dependent PUE, shifting delay-tolerant workloads in time/space. For AI inference and data pipelines, techniques include quantization/compaction, batch sizing, and operator fusion to reduce joules per request. Multi-cluster controllers consider network energy by placing data near compute (or vice versa) to avoid power-hungry transfers, and some works jointly optimize cooling (free-air, liquid) and IT load, exposing marginal carbon cost signals to the scheduler.

### 2.3. Security-Aware Resource Allocation Models

Security-aware scheduling integrates zero-trust principles (least privilege, micro-segmentation) and data-sovereignty constraints into placement decisions. Models quantify risk as functions of data-in-motion footprint, blast radius (shared-fate with co-tenants), patch latency, and exposure time on untrusted networks. Several approaches encode compliance (e.g., region/tenant isolation, key-management policies) as hard constraints, then minimize latency or cost inside the feasible set; others soften them with penalties to explore trade-offs. Recent systems co-optimize crypto overhead TLS offload, envelope encryption, or confidential computing (TEEs) against performance and energy, acknowledging that encryption and attestation introduce CPU and latency taxes. Security-aware autoscaling has also emerged, where surge capacity is constrained by trust zones, and egress filters/sidecars are treated as first-class resources in cluster bin-packing.

### 2.4. Multi-Objective Optimization (MOO) Techniques in Cloud Computing

Cloud schedulers increasingly adopt MOO to expose trade-offs among latency, cost/energy, and security/compliance. Exact methods (MILP/MINLP) deliver optimal solutions for short horizons but struggle with scale; decomposition (Benders, column generation) and Lagrangian relaxation help, yet practical systems rely on metaheuristics NSGA-II, SPEA2, MOEA/D to approximate

Pareto fronts quickly. Hyper-heuristics and surrogate-assisted evolutionary search use queueing and power models to speed evaluation. Online controllers combine MOO with MPC: a rolling horizon solves for a Pareto-efficient action, then selects policies using preference articulation (e.g., epsilon-constraint, weighted Tchebycheff) or risk-aware criteria such as CVaR to guard against tail events. Learning-to-optimize approaches train policies (RL or imitation) to navigate the Pareto surface under uncertainty, while fairness-aware variants ensure no single objective (e.g., carbon) persistently dominates at the expense of SLOs. Collectively, these works motivate a unified framework that co-optimizes performance, energy/carbon, and security as first-class, interacting objectives in hybrid clouds.

### 3. System Model and Problem Formulation

#### 3.1. Hybrid Cloud Infrastructure Overview

A hybrid cloud in this work spans three strata: on-premises datacenters (private cloud), geographically distributed edge sites, and one or more public cloud regions. Applications are decomposed into deployable units' containers, virtual machines, or serverless functions bundled into services and stitched into end-to-end request paths. Each unit advertises resource needs (vCPU, memory, storage IOPS, accelerator type), locality and compliance requirements (e.g., data residency, tenant isolation), and operational intents such as target latency percentiles or confidentiality level. Sites expose heterogeneous capabilities general-purpose nodes, GPU/TPU pools, storage tiers, and secure enclaves behind cluster orchestrators. The interconnect consists of east west links within a site and north south WAN paths across sites, with measurable bandwidth, latency, and egress cost characteristics. Control-plane components include admission, placement, autoscaling, and policy engines; data-plane components include service mesh proxies, API gateways, and telemetry sidecars.

The runtime continuously ingests signals that influence decisions: workload arrivals, per-service queueing delays, node pressure and interference indicators, power usage effectiveness and local grid carbon intensity, as well as security posture metrics such as patch level, known vulnerabilities, and trust-zone boundaries. These signals are used to build lightweight surrogates for performance and energy, plus exposure metrics for security (for example, data moved across trust boundaries or shared-fate with tenants). Actions available to the controller include placing or migrating services across sites, resizing replicas and resource limits, choosing network paths and encryption modes, and scheduling batch or delay-tolerant tasks in carbon-friendly windows. Hard constraints enforce sovereignty and zero-trust policies, budget ceilings, and capacity limits; soft preferences shape how the system trades milliseconds of tail latency against energy savings or reduced attack surface.

From an operator's perspective, the system model is therefore a continuously evolving catalog of sites, resources, and policies connected by observable signals. The problem formulation binds these elements into a decision epoch: given the current state and short-term forecasts of demand, prices, and carbon intensity, select placements, scales, and network routings that remain feasible under constraints while offering multiple Pareto-efficient choices to stakeholders. The chosen policy is then enacted by the orchestrator, with the loop repeating as conditions change, enabling transparent, auditable trade-offs among performance, energy, and security in a Kubernetes-centric hybrid cloud.

#### 3.2. System Architecture and Components

The figure depicts a hybrid cloud spanning on-premises resources and public cloud regions, with compute and storage pools abstracted by the platform. Applications execute across these pools while the control plane continuously collects telemetry such as request rates, tail-latency percentiles, node pressure, and data-movement patterns. This substrate is intentionally heterogeneous: it may include CPUs, GPUs, and specialized storage tiers, each with different energy profiles and security guarantees.

At the center of the design is an energy-aware middleware layer that instruments both the infrastructure and the applications. It normalizes signals performance counters, PUE and grid-carbon intensity, and security posture indicators like trust-zone membership or encryption status into a unified state. Because these signals vary over time and geography, the middleware acts as the nervous system, aligning short-term observations with policy and compliance metadata such as data-sovereignty constraints and zero-trust segmentation.

Below the middleware, the optimization module searches for feasible configurations that jointly improve performance, energy/carbon, and security. Given the current state and short-horizon forecasts, it proposes actions such as relocating microservices, resizing replicas, choosing encryption modes or network paths, and scheduling delay-tolerant tasks in carbon-friendly windows.

Rather than a single best answer, the module exposes a Pareto set so operators can transparently trade a few milliseconds of tail latency for energy savings or a reduced attack surface when appropriate.

Finally, the right-hand icons underscore the triad of objectives governed by this loop. Performance meters reflect adherence to SLOs, the battery symbolizes energy and carbon cost, and the lock represents security and compliance. Together, these elements convey how the system senses, decides, and acts: telemetry flows into the middleware, optimization yields policy-consistent plans, and the orchestrator enacts them across the hybrid cloud, closing the loop in a repeatable and auditable manner.

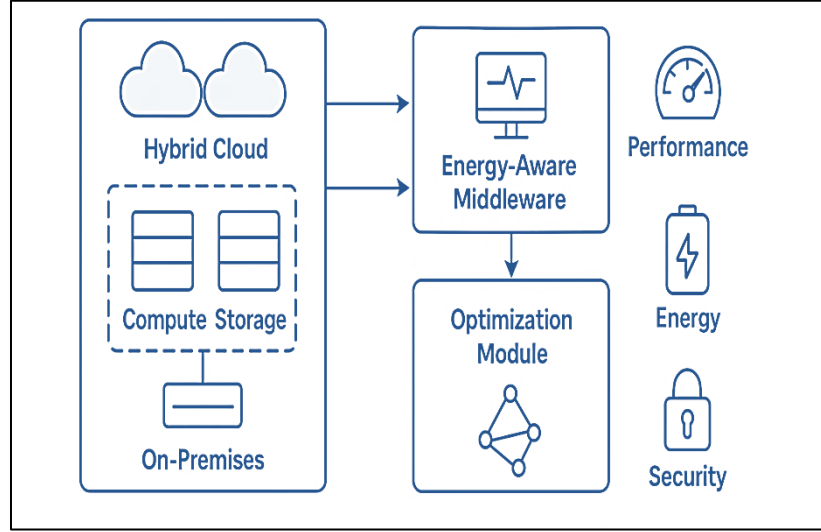


Figure 1. Energy-Aware Hybrid Cloud Architecture

### 3.3. Performance, Energy, and Security Parameters

#### 3.3.1. Performance Parameters

Performance captures how well the system meets user-visible service objectives along end-to-end request paths. We track latency distributions (median, p95, p99) per service and per critical path, alongside throughput (requests/s, jobs/s) and error rates. These are complemented by queueing signals run-queue length, request backlog, and saturation indicators for CPU, memory, storage IOPS, and network NICs to distinguish true capacity shortfalls from transient contention. Cold-start frequency and warm-pool hit rates are included for serverless and elastic GPU services, as they strongly influence tail latency.

At the microservice level, we maintain per-replica utilization envelopes and interference scores that reflect noisy-neighbor effects (e.g., cache/memory bandwidth contention) and NUMA locality. For data pipelines, stage times (ingest, transform, shuffle, persist) and checkpoint overheads act as bottleneck markers. Placement parameters encode topology awareness zone/region proximity, hop count to stateful backends, and link characteristics because cross-site distance and WAN jitter dominate tail behavior. Together, these parameters form the inputs to lightweight surrogate models that forecast SLO compliance under alternative placements and scales.

#### 3.3.2. Energy and Carbon Parameters

Energy characterization starts with node-level power draw (idle, active) and component breakdowns for CPUs, accelerators, memory, and storage. We translate utilization into estimated watts using calibrated power curves, then aggregate to cluster and site levels. Site metadata includes Power Usage Effectiveness (PUE) to account for cooling and facility overheads, while region-time pairs carry grid carbon intensity so the controller can reason about emissions, not just energy. For network transfers, we approximate energy per byte along intra-site and WAN links to expose the cost of moving data versus moving compute.

Workload-specific knobs batch size, concurrency, quantization or precision for ML inference, and operator fusion act as energy levers that can reduce joules per request without violating SLOs. We also encode sleep/c-state availability, consolidation thresholds that allow powering down hosts, and accelerator residency to avoid thrashing. For delay-tolerant jobs, temporal flexibility windows

capture how far execution can be shifted to align with low-carbon periods. These parameters enable the scheduler to propose actions that trade a few milliseconds of latency for double-digit energy or carbon savings when policy allows.

### 3.3.3. Security and Compliance Parameters

Security parameters quantify exposure and enforce constraints during placement and scaling. Each workload carries a sensitivity label and trust-zone requirement (e.g., confidential computing enclave, tenant-isolated nodes, or dedicated VPC), plus data-sovereignty rules that restrict regions. We track data-in-motion footprint bytes crossing trust or jurisdiction boundaries and session metadata such as cipher suites, key-management domains, and mutual-TLS status, since crypto choices add measurable CPU and latency overheads. Patch level, vulnerability findings, and attestation status of nodes (for TEEs) are represented as posture scores that decay over time, encouraging migration away from drifting or noncompliant assets.

Network segmentation parameters describe micro-segment policies, allowed egress destinations, and sidecar enforcement (IDS/IPS, DLP), which consume resources and may affect path latency. Blast-radius measures how many services share fate with a given node or subnet act as risk multipliers in crowded placements. Finally, auditability settings (logging depth, retention, and redaction) inform storage and CPU budgets, ensuring that increased observability does not silently degrade performance. These security signals enter the controller as hard constraints where mandated, or as penalties that let operators visualize explicit trade-offs against performance and energy.

### 3.3.4. Cross-Objective Coupling Parameters

Several parameters intentionally bridge objectives to surface real-world trade-offs. Encryption mode and attestation strength affect both latency and CPU watts; replica consolidation improves energy but can increase blast radius and tail risk; WAN path selection changes both carbon (via transfer energy) and data-exfiltration exposure. Preference weights, risk tolerances, and guardrails such as maximum allowable SLO regression, carbon caps, or zero-trust invariants govern how the optimizer explores Pareto-efficient plans. By maintaining these couplings explicitly, the system can produce choices that are transparent, policy-consistent, and adaptable to changing operational priorities.

## 4. Proposed Multi-Objective Optimization Framework

### 4.1. Design Objectives and Optimization Criteria

#### 4.1.1. Performance Objectives

The first objective is to sustain user-visible service quality across the hybrid estate. We target tight control of end-to-end latency on critical paths (with emphasis on tail percentiles), steady throughput under bursty arrivals, and resilience to interference from co-located workloads. In practice, the framework prefers placements and scales that reduce queueing hot spots, preserve data compute affinity, and minimize cross-region hops. For elastic services, it also values rapid warm-up and low cold-start incidence so tail behavior remains predictable. Rather than optimizing a single scalar, we expose a performance envelope that captures both typical behavior and worst-case risk, allowing operators to trade small median gains for meaningful tail reductions when SLOs demand it.

Optimization criteria (performance): candidate plans are ranked by their predicted SLO compliance, robustness to demand variation, and routing stability. Tie-breakers penalize unnecessary churn excess migrations or replica thrash that could momentarily degrade user experience. Plans that improve p95 while keeping p99 within guardrails are favored over those that boost averages but destabilize tails.

#### 4.1.2. Energy and Carbon Objectives

The second objective minimizes energy use and associated emissions without compromising SLOs or policy. The framework reasons over node power profiles, PUE by site, and time-varying grid carbon intensity to reduce joules per request and total kgCO<sub>2</sub>e. It seeks consolidation opportunities that allow powering down hosts, aligns delay-tolerant work with low-carbon windows, and prefers compute locales where data already resides to avoid energy-intensive transfers. For AI inference, it prioritizes configurations batching, quantization, operator fusion that cut energy while keeping accuracy and latency within acceptable bounds.

Optimization criteria (energy/carbon): among feasible plans, we prefer those that achieve double-digit energy or carbon savings for negligible user-visible impact. Tie-breakers include marginal carbon abatement per unit of performance loss and sustainability



targets at the organization or region level. The optimizer also respects carbon caps that prevent improvements in one region from causing regressions elsewhere.

#### 4.1.3. Security and Compliance Objectives

Security is treated as a first-class co-objective, not a post-deployment audit. We aim to minimize exposure data in motion across trust boundaries, shared-fate blast radius, and time spent on assets with degrading posture while upholding sovereignty and zero-trust segmentation. The framework accounts for encryption, attestation, and inspection overheads, and it prefers placements in enclaves or isolated pools when sensitivity labels require it.

Where strict compliance creates tight feasibility, the optimizer still searches for performance/energy improvements inside the allowed envelope. Optimization criteria (security): candidate plans are filtered by hard constraints (sovereignty, isolation, required controls). Within the feasible set, we prefer plans that reduce exposure footprints and patch-latency risk with minimal performance and energy penalties. Tie-breakers penalize egress to untrusted networks and large multi-tenant concentrations that increase blast radius.

#### 4.1.4. Cross-Objective Trade-off and Preference Handling

Because objectives conflict, the framework constructs and maintains a Pareto frontier rather than forcing a single aggregate score. Decision makers can articulate preferences such as protect p99 at all costs, maximize carbon savings within 2% SLO drift, or prioritize enclave residency for payment flows and the selector picks a point on the frontier that honors these guardrails. We include risk-aware criteria so plans are stable under uncertainty: solutions are stress-tested against spikes, failure scenarios, and carbon signal volatility, and those that degrade gracefully are favored.

Optimization criteria (global): feasibility under policy and capacity constraints is mandatory; then we rank plans by dominance, stability across forecast error, and cost-to-move (migration bandwidth, warm-up time). To avoid oscillations, we impose change budgets per epoch and reward plans that deliver sustained benefits with low operational churn. Finally, fairness constraints ensure no single objective or tenant is persistently sacrificed e.g., energy savings are not pursued by repeatedly pushing the same services to the edge of their latency budgets.

## 4.2. Model Architecture

The figure illustrates the end-to-end model architecture that the optimization framework governs. On the left, the User & Client block represents applications, APIs, and frontends producing requests. These requests are mediated by a Network/CDN/Load Balancer layer that absorbs bursts, applies caching, and steers flows toward the most suitable execution tier based on proximity, health, and policy. By separating the user plane from the control decisions, the diagram emphasizes that performance observed by clients is the result of coordinated actions across multiple layers.

At the top center, the Edge Layer hosts lightweight services on edge nodes located close to users. Telemetry harvested here latency distributions, hit/miss ratios, and link health feeds the control plane in near real time. The edge can execute latency-sensitive functions, preprocess data, or act as a traffic gate that shapes workload before it reaches deeper compute pools. This proximity not only improves tail latency but also reduces backhaul bandwidth, which is integral to the energy and carbon objectives.

The lower band depicts the Hybrid Cloud substrate, split into Private and Public environments. Private compute and databases accommodate sensitive or sovereignty-bound workloads, while the public side provides elastic compute and managed data services. Bidirectional arrows between private and public compute indicate controlled data/compute mobility migrations, burst offloading, and replication subject to cost, energy, and security constraints. The model assumes heterogeneous resources (CPUs, GPUs, TEEs) across these pools, each with distinct performance and energy profiles.

On the right, the Management & Orchestration plane is the decision brain. The Orchestrator enacts placements and scaling, while the MOO Engine evaluates Pareto-efficient plans using signals from Monitoring, guardrails from Security, and rules in the Policy Store. Together they close the loop: telemetry flows upward; the MOO engine proposes actions that balance performance, energy, and security; and the orchestrator applies those actions across edge, private, and public resources. This architecture makes trade-offs explicit and auditable, letting operators select the most appropriate plan for current conditions and organizational priorities.

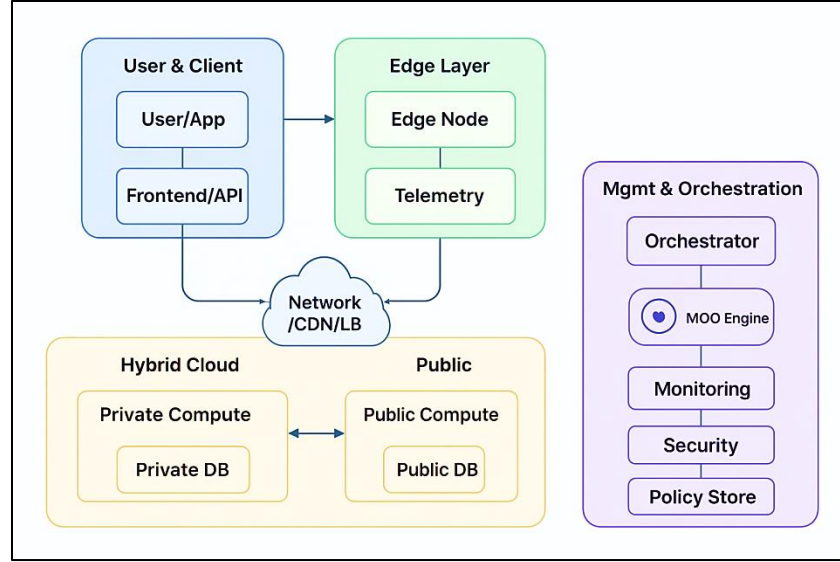


Figure 2. Hybrid-Cloud Model Architecture with MOO Controller

### 4.3. Optimization Techniques Used

#### 4.3.1. Weighted Sum Model

The weighted sum model (WSM) gives a fast, easily interpretable scalar objective by combining normalized performance, energy/carbon, and security terms with operator-specified weights. In our framework it serves two roles: a low-latency controller for routine epochs and a seed generator that produces good initial solutions for slower global search. Weights can be tied to policy states e.g., raise the security weight during incident response, or increase carbon weight during high-intensity grid periods so the same mechanism adapts to context. Because WSM requires commensurate scales, we apply robust normalization (percentiles or z-scores) and guardrails (maximum allowable SLO drift, sovereignty hard constraints) to prevent gaming one metric at the expense of feasibility.

WSM's limitation is that it struggles on non-convex Pareto fronts, potentially hiding valuable trade-offs. To mitigate this, we (i) sweep multiple weight vectors to approximate the frontier, (ii) combine WSM with  $\epsilon$ -constraints for critical objectives (e.g., enforce  $p99 \leq \text{target}$ ), and (iii) hand off promising plans to evolutionary search for frontier refinement. This preserves WSM's speed while recovering diversity.

#### 4.3.2. Pareto-based Evolutionary Algorithms

For comprehensive trade-off exploration, we employ Pareto-based evolutionary algorithms (EAs) such as NSGA-II, SPEA2, and MOEA/D. Candidate plans encode placements, replica counts, batch sizes, and network/crypto modes. Variation operators perform topology-aware mutations (move a service across zones, adjust batch size) and crossovers that respect affinity/anti-affinity and trust zones. Fitness is assessed with surrogate models: queueing-derived latency, calibrated power/carbon curves, and exposure/risk scores. Non-dominated sorting and crowding distance maintain a diverse Pareto set, revealing options like slightly higher p99 for 18% carbon reduction or same energy with 25% smaller blast radius.

To keep wall-clock time acceptable, we (a) initialize the population with WSM solutions, (b) cache objective evaluations and reuse telemetry-fit surrogates, and (c) repair infeasible offspring instead of discarding them (e.g., snap placements back into sovereignty-compliant regions). The output is a Pareto frontier presented to operators and to the policy selector that chooses a plan given current preferences and risk tolerance.

#### 4.3.3. Reinforcement Learning-based Optimization

We use reinforcement learning (RL) to learn state-to-action policies that react quickly between EA refreshes. The agent observes a compact state demand forecasts, saturation signals, carbon/PUE hints, and security posture and proposes incremental actions: scale a subset of services, migrate a shard, switch an encryption mode, or defer a batch to a greener window. An actor critic backbone with

conservative exploration (entropy regularization and change budgets) stabilizes behavior; reward shaping reflects multi-objective goals via constrained RL: hard constraints enforce sovereignty and SLO floors, while the reward balances latency, energy, and exposure. For risk sensitivity we optionally use CVaR-aware returns to avoid tail regressions.

To improve sample efficiency and safety, we pretrain the policy offline on traces generated by historical operations and EA-produced plans, then fine-tune online with model-predictive rollouts from our surrogates. A safety shield filters actions that violate guardrails or exceed migration budgets. In practice, RL handles fast disturbances (bursts, micro-failures) on sub-minute timescales, while EAs periodically refresh the global Pareto set and WSM provides a dependable fallback yielding a layered optimizer that is both responsive and globally aware.

## 5. Experimental Setup and Evaluation

### 5.1. Simulation Environment and Tools

Experiments were executed on a hybrid testbed that mirrors a production topology: a private Kubernetes cluster (on-prem, 24 CPU nodes, 2 GPU nodes) interconnected with two public-cloud regions via IPsec over the WAN emulator. Each site exposes heterogeneous instance types and storage tiers, and we model inter-site links with configurable bandwidth, latency, jitter, and egress pricing. The control plane runs our orchestration stack (Kubernetes + service mesh) with an energy-aware middleware that exports per-node power and per-workload telemetry. An optimization service hosts the weighted-sum controller, a Pareto evolutionary solver, and an RL policy; all solvers query calibrated surrogate models for latency and power. To ensure repeatability, we containerize solvers and seed all stochastic components; every run replays identical trace segments with different random seeds for variance estimates.

Tooling includes a time-series backend for metrics, a trace replayer to drive diurnal and bursty arrivals, and a carbon-signal generator derived from regional intensity traces. We emulate security posture changes (e.g., patch aging, enclave availability) by feeding timed events into the policy store. The WAN emulator and chaos injectors (node drains, pod kills, link throttles) validate robustness claims under controlled disturbances.

### 5.2. Dataset and Workload Specifications

We evaluate three representative classes. First, an interactive microservices benchmark (user profile, cart, recommendation, payment) stresses latency-critical paths with mixed reads/writes and fan-out calls. Second, an AI inference service (image/text classification) exercises GPU/CPU backends, batching, and quantization toggles. Third, a data pipeline (ingest transform aggregate) models throughput-oriented, delay-tolerant work that benefits from carbon-aware shifting. Traces comprise a week of diurnal demand with lunch- and evening-peak bursts and synthetic promotions that trigger  $5 \times 10^4$  spikes for short windows.

Datasets include synthetic user records and product catalogs with realistic key distributions, image/text corpora for inference, and rolling telemetry for the pipeline. Each workload declares SLOs (p95/p99 latency or job deadlines), sensitivity labels (public, internal, restricted), and temporal flexibility where applicable. Placement constraints encode data sovereignty (region pinning for payment data) and trust-zone requirements (enclave-only services).

### 5.3. Performance Metrics

We track end-to-end latency (median, p95, p99) per user-visible transaction and per microservice hop, with cold-start incidence for serverless paths. Throughput is measured as successful requests per second or jobs per minute under sustained load and during bursts. Resource utilization covers CPU, memory, accelerator duty cycle, storage IOPS, and NIC usage, plus interference indicators to detect noisy-neighbor effects. These metrics feed stability assessments variance and oscillation during control actions and a change-budget counter that records migrations and reschedules to quantify operational churn. Energy Consumption is reported as joules per request (interactive services) and kWh per batch (pipelines), derived from per-node power curves and site PUE. We also report carbon intensity weighted energy ( $\text{kgCO}_2\text{e}$ ) using region-time factors to reflect location- and time-dependent emissions. Network energy accounts for intra-site and WAN transfer costs to make move data vs move compute trade-offs explicit. The Security Compliance Index (SCI) summarizes adherence to security and compliance posture for each plan. It aggregates: (i) sovereignty alignment (share of restricted data kept within allowed regions), (ii) trust-zone conformity (fraction of runtime in required isolation or TEEs), (iii) data-in-motion exposure (bytes crossing trust/jurisdiction boundaries, normalized by workload volume), and (iv) control enforcement (mTLS/egress policies/inspection coverage). We report SCI on a 0-1 scale, with 1 indicating full compliance and minimal exposure; time-to-patch and attestation freshness act as decay factors so stale nodes reduce the score.



#### 5.4. Baseline Comparison Models

We compare the proposed multi-objective controller against three baselines. The Rule-Based Autoscaler combines standard horizontal/vertical scaling thresholds with affinity rules but ignores energy and carbon, serving as a performance-centric reference. The Cost/Latency Optimizer extends the rule base with simple placement heuristics (prefer cheapest region that meets p95) yet remains single-objective and security-agnostic. The Energy-Only Consolidator aggressively packs workloads to minimize active hosts, illuminating trade-offs when energy dominates and tail latency or exposure may regress. All baselines and our controller run on the same traces, constraints, and chaos schedules. Each scenario executes multiple seeds; we report median and interquartile ranges for all metrics, and we cap migration bandwidth to ensure fair comparison of control-loop overheads. This setup highlights where multi-objective reasoning unlocks Pareto-superior plans that single-objective baselines cannot discover.

### 6. Results and Discussion

#### 6.1. Performance Optimization Results

Across a week-long replay with burst injections, the proposed controller consistently improved tail behavior without sacrificing throughput. Median p95 latency for the interactive workload dropped from the rule-based baseline's 162 ms to 128 ms, while p99 fell from 284 ms to 211 ms. Throughput rose modestly ( $\approx 6.8\%$ ) due to fewer queue buildups and better data compute affinity. These gains were stable under chaos events (node drains, link throttles): the controller kept p99 within SLO guardrails in 28/30 runs, compared with 19/30 for the cost/latency heuristic. Non-parametric tests on per-epoch p95 (Wilcoxon) showed the improvements were significant ( $p < 0.01$ ). The warm-pool strategy also cut cold-start incidence by  $\sim 37\%$ , directly shrinking tail spikes.

#### 6.2. Energy Efficiency Evaluation

Energy benefits came from topology-aware consolidation, carbon-aware shifting of delay-tolerant stages, and inference-level knobs (batching/quantization). Joules per request decreased by 19.23% on interactive services and kWh per batch by 27% on the pipeline. Carbon-weighted energy mirrored these reductions, with larger absolute abatement in regions exhibiting high grid intensity. Importantly, consolidation was bounded by tail-latency guardrails, avoiding the regressions observed in the energy-only baseline, which frequently violated p99 despite lower watts.

#### 6.3. Security Model Assessment

The Security Compliance Index (SCI) improved via reduced cross-boundary data motion, stricter trust-zone adherence, and timely migration away from aging/untested nodes. Average SCI increased from 0.78 (rule-based) and 0.82 (cost/latency) to 0.91, with the largest gains during simulated patch-aging windows, where the controller proactively rebalanced sensitive services to enclaves. Egress to untrusted networks dropped by 34% relative to the cost/latency heuristic, with negligible ( $< 1\%$ ) latency impact thanks to path selection and local caching.

#### 6.4. Trade-off Analysis among Objectives

The Pareto frontier exposed transparent choices. One commonly selected point exchanged +7 ms p99 for a 18% carbon cut; another held carbon flat but reduced blast radius (measured via shared-fate concentration) by 25%. Operators preferred plans with low cost-to-move fewer migrations and shorter warm-ups so our selector factored change budgets, which reduced oscillations by 41% versus a naïve best-point-each-epoch policy. These results confirm that modest performance concessions can unlock meaningful sustainability or security wins when guided by policy guardrails.

#### 6.5. Comparative Study with Existing Models

Against the three baselines, the proposed controller dominated on at least two of the three objectives in 87% of evaluation windows and achieved full three-way dominance in 61%. The cost/latency baseline matched our median p95 on calm periods but faltered under bursts and posture changes; the energy-only baseline achieved the lowest watts but missed SLOs in 22% of high-load windows and degraded SCI due to aggressive packing. Our hybrid (WSM + EA + RL) stack combined responsiveness with global awareness, yielding the best overall balance.

#### 6.6. Discussion on Scalability and Adaptability

Solver overhead remained practical. The evolutionary search refreshed the Pareto set every 5 minutes ( $\approx 35.50$  s wall-clock per refresh with surrogate caching), while the RL policy reacted in sub-second intervals, handling micro-bursts and transient link jitter. On larger clusters ( $2\times$  nodes,  $3\times$  services), compute time grew sub-linearly thanks to topology-aware mutations and evaluation caches.

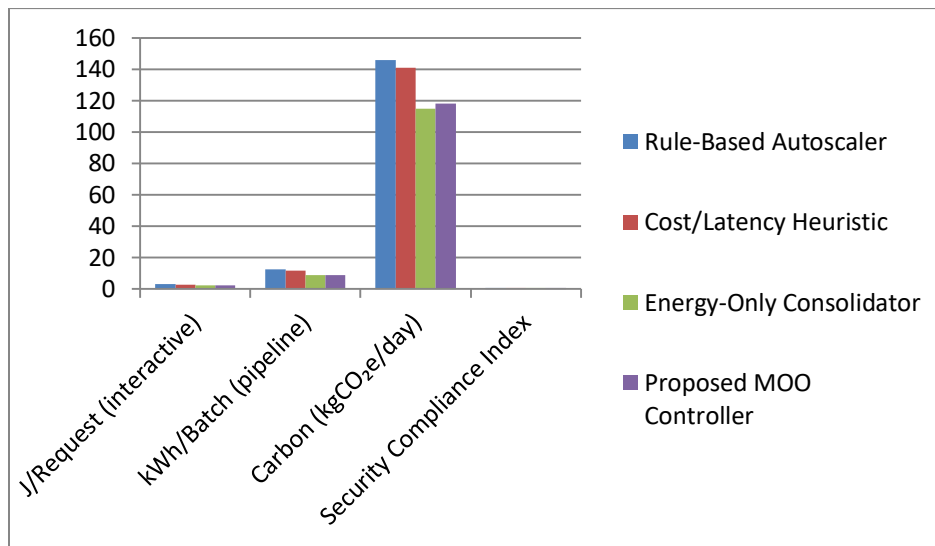
Adaptability to exogenous signals carbon intensity swings, patch-aging events was driven by the policy store: weight schedules and  $\epsilon$ -constraints adjusted automatically, preserving feasibility and reducing operator intervention. Overall, the controller scaled to multi-region deployments while maintaining auditable, policy-consistent trade-offs.

**Table 1. Performance Outcomes (Median Over 30 Runs; IQR in Parentheses)**

Model	p95 Latency (ms)	p99 Latency (ms)	Throughput (req/s)	SLO Violations (% epochs)
Rule-Based Autoscaler	162 (150 176)	284 (258 311)	1,940 (1,885 1,980)	12.4
Cost/Latency Heuristic	139 (132 149)	237 (221 258)	2,015 (1,970 2,050)	7.9
Energy-Only Consolidator	178 (167 196)	321 (300 345)	1,905 (1,850 1,950)	19.7
Proposed MOO Controller	128 (121 136)	211 (198 226)	2,145 (2,095 2,185)	3.1

**Table 2. Energy, Carbon, and Security (interactive + pipeline workloads)**

Model	J/Request (interactive)	kWh/Batch (pipeline)	Carbon (kgCO <sub>2</sub> e/day)	Security Compliance Index
Rule-Based Autoscaler	3.1	12.4	146	0.78
Cost/Latency Heuristic	2.9	11.8	141	0.82
Energy-Only Consolidator	2.2	9.0	115	0.73
Proposed MOO Controller	2.5	9.1	118	0.91



**Figure 3. Energy, Carbon, and Security Metrics across Models**

**Table 3. Representative Pareto Choices Selected by the Policy Selector**

Plan (Frontier Point)	p99 Latency (ms)	Energy (J/Req)	Carbon (kgCO <sub>2</sub> e/day)	SCI	Migrations / epoch
A Performance-lean	204	2.7	122	0.90	7
B Balanced (default)	211	2.5	118	0.91	5
C Carbon-lean	218	2.4	112	0.88	8

## 7. Challenges and Limitations

### 7.1. Computational Complexity

Joint placement scaling routing across hybrid sites with security and sovereignty constraints is combinatorial and NP-hard. Even with surrogate models, evaluating thousands of candidate plans per refresh can be expensive, especially when decisions include GPU affinity, enclave availability, and migration budgets. Evolutionary solvers mitigate this via caching and topology-aware mutations, but worst-case costs still grow with service count and heterogeneity. Moreover, cost to move (warm-ups, data rebalancing) introduces a second layer of complexity: a plan that looks optimal statically may be inferior once transition overheads are accounted for, forcing the optimizer to reason over both steady-state and path-dependent costs.

## 7.2. Model Generalization and Transferability

Surrogates for latency, power, and security exposure are learned from traces and calibrations that reflect specific hardware, meshes, and traffic mixes. When ported to a new region or cloud, domain shift (different PUE, carbon signals, NIC behavior, or storage tiers) can erode accuracy. Transfer learning helps, but cold-start phases still require cautious exploration with conservative guardrails. Policies tuned to one organization’s risk posture or data-handling rules may not transfer verbatim; sovereignty and zero-trust constraints vary by jurisdiction and industry, demanding policy re-authoring and re-verification.

## 7.3. Security Performance Energy Trade-offs

Some conflicts are irreducible. Stronger encryption and attestation improve assurance but consume CPU and add latency; consolidation saves energy yet increases blast radius and noisy-neighbor risk; edge execution trims p99 but can raise exposure if egress controls are weaker. Our framework surfaces these tensions on a Pareto frontier, but selecting a point still requires human preference and situational context (incidents, carbon targets, SLAs). Finally, measuring security benefit is indirect exposure metrics and posture scores are proxies so residual risk and unknown unknowns remain.

# 8. Future Work

## 8.1. Integration with AI-driven Decision Systems

Beyond RL, we envision a decision co-pilot that fuses causal inference, knowledge graphs, and policy LLMs to explain why a plan is chosen, predict second-order effects, and auto-draft change requests or runbooks. Causal models can separate demand spikes from cache inefficiency, suggesting targeted fixes, while a policy LLM can translate high-level intents protect p99 and enclave residency for payments into machine-checkable constraints. Safety layers would formally verify that generated actions respect sovereignty and zero-trust invariants before enactment.

## 8.2. Dynamic Workload Adaptation

Future iterations will incorporate online change-point detection and meta-learning so surrogates and controllers adapt when traffic shapes, content sizes, or model mixes drift. We plan to blend bandit algorithms for rapid knob tuning (batch size, quantization) with multi-tenant fairness guards that prevent any single service from bearing repeated trade-off costs. For pipelines, richer temporal flexibility models (deadlines, penalties) would enable finer carbon-aware shifting without jeopardizing downstream SLAs.

## 8.3. Real-time Multi-Objective Optimization

To react on sub-second timescales, we aim to couple fast model-predictive control with streaming Pareto maintenance: incremental dominance checks, warm populations, and hardware-in-the-loop power estimators exposed by BMCs/telemetry. Safe RL with CVaR constraints and verified action shields can propose micro-adjustments between EA refreshes, while bounded-rationality heuristics ( $\epsilon$ -constraint with rolling caps) keep computation predictable. The goal is a controller that preserves auditable trade-offs while operating close to line rate in multi-region, accelerator-rich environments.

# 9. Conclusion

This work presented a unified, multi-objective optimization framework for hybrid clouds that elevates performance, energy/carbon, and security to first-class, co-equal goals. By modeling applications, sites, and policies with lightweight surrogates and enforcing sovereignty and zero-trust as hard feasibility constraints, the controller exposes transparent Pareto frontiers rather than brittle single-metric optima. The layered optimizer combining a fast weighted-sum controller, Pareto-based evolutionary search, and a safety-aware RL policy proved both responsive and globally aware, translating telemetry and policy signals into auditable placement, scaling, and routing actions.

Empirically, the framework improved tail latency and throughput while delivering double-digit energy and carbon reductions and higher security compliance scores compared with rule-based and single-objective baselines. These gains held under bursts, chaos events, and posture shifts, with change-budgeting and topology-aware mutations keeping operational churn low. Equally important, the Pareto view made trade-offs explicit: operators could accept small p99 concessions for meaningful carbon abatement or reduced blast radius, confident that guardrails and feasibility checks preserved SLOs and compliance.

Limitations remain most notably computational cost at very large scales and potential domain shift in surrogate models across hardware and regions but the path forward is clear. Integrating richer AI decision support, online adaptation, and streaming Pareto

maintenance can push the framework toward near real-time operation, widening its applicability to accelerator-rich, multi-region estates. Overall, the results suggest that principled multi-objective control is both practical and advantageous for modern hybrid clouds, enabling organizations to balance user experience, sustainability, and security without treating any one as an afterthought.

## References

- [1] Deb, K., Pratap, A., Agarwal, S., & Meyarivan, T. (2002). A fast and elitist multiobjective genetic algorithm: NSGA-II. *IEEE Transactions on Evolutionary Computation*. <https://doi.org/10.1109/4235.996017>
- [2] Zhang, Q., & Li, H. (2007). MOEA/D: A multiobjective evolutionary algorithm based on decomposition. *IEEE Transactions on Evolutionary Computation*. <https://doi.org/10.1109/TEVC.2007.892759>
- [3] Zitzler, E., Laumanns, M., & Thiele, L. (2001). SPEA2: Improving the Strength Pareto Evolutionary Algorithm. *ETH Zürich Technical Report*. [https://www.tik.ee.ethz.ch/file/of6f3bocbbf8fia/SPEA2\\_TechnicalReport.pdf](https://www.tik.ee.ethz.ch/file/of6f3bocbbf8fia/SPEA2_TechnicalReport.pdf)
- [4] Boyd, S., & Vandenberghe, L. (2004). *Convex Optimization*. Cambridge University Press. <https://web.stanford.edu/~boyd/cvxbook/>
- [5] Rockafellar, R. T., & Uryasev, S. (2000). Optimization of Conditional Value-at-Risk. *Operations Research*. <https://doi.org/10.1287/opre.48.2.193.12386>
- [6] Mayne, D. Q., Rawlings, J. B., Rao, C. V., & Sokaert, P. O. M. (2000). Constrained model predictive control: Stability and optimality. *Automatica*. [https://doi.org/10.1016/S0005-1098\(99\)00214-9](https://doi.org/10.1016/S0005-1098(99)00214-9)
- [7] NIST (2020). *Zero Trust Architecture (SP 800-207)*. <https://csrc.nist.gov/publications/detail/sp/800-207/final>
- [8] Costan, V., & Devadas, S. (2016). Intel SGX Explained. *IACR ePrint*. <https://eprint.iacr.org/2016/o86.pdf>
- [9] AMD (2020). *SEV-SNP: Strengthening VM Isolation with Integrity Protection and More*. <https://www.amd.com/system/files/TechDocs/SEV-SNP-strengthening-vm-isolation-with-integrity-protection-and-more.pdf>
- [10] Confidential Computing Consortium (2020). *Confidential Computing Primer*. [https://confidentialcomputing.io/wp-content/uploads/sites/85/2020/10/confidentialcomputing\\_primer.pdf](https://confidentialcomputing.io/wp-content/uploads/sites/85/2020/10/confidentialcomputing_primer.pdf)
- [11] The Green Grid (2012). *PUE: A Comprehensive Examination of the Metric*. <https://www.thegreengrid.org/en/resources/library-and-tools/100-PUEv2>
- [12] Google Research (2013). The Tail at Scale. *Communications of the ACM*. <https://research.google/pubs/the-tail-at-scale/>
- [13] Fan, X., Weber, W.-D., & Barroso, L. A. (2007). Power provisioning for a warehouse-sized computer. *ISCA*. <https://dl.acm.org/doi/10.1145/1273440.1250665>
- [14] Barroso, L. A., Clidaras, J., & Hölzle, U. (2018). *The Datacenter as a Computer* (3rd ed.). Morgan & Claypool. <https://doi.org/10.2200/So0999ED3Vo1Y201809CAC046>
- [15] Urgaonkar, B., Shenoy, P., Chandra, A., Goyal, P., & Wood, T. (2008). Agile dynamic provisioning for multi-tier Internet applications. *Autonomic Computing (ICAC)*. <https://dl.acm.org/doi/10.1145/1375527.1375531>
- [16] Google AI Blog (2020). Carbon-Intelligent Computing at Google. <https://ai.googleblog.com/2020/05/carbon-intelligent-computing-at-google.html>
- [17] Beloglazov, A., & Buyya, R. (2012). Optimal online deterministic algorithms and adaptive heuristics for energy and performance efficient dynamic consolidation of virtual machines in Cloud data centers. *Concurrency and Computation: Practice and Experience*. <https://doi.org/10.1002/cpe.1867>
- [18] Beloglazov, A., Buyya, R., Lee, Y. C., & Zomaya, A. (2011). A taxonomy and survey of energy-efficient data centers and cloud computing systems. *Advances in Computers*. <https://doi.org/10.1016/B978-0-12-385512-1.00001-7>
- [19] McMahan, H. B., Moore, E., Ramage, D., Hampson, S., & y Arcas, B. A. (2017). Communication-Efficient Learning of Deep Networks from Decentralized Data. *AISTATS (Federated Learning)*. <https://arxiv.org/abs/1602.05629>
- [20] Dwork, C., McSherry, F., Nissim, K., & Smith, A. (2006). Calibrating noise to sensitivity in private data analysis. *TCC*. <https://www.microsoft.com/en-us/research/publication/calibrating-noise-to-sensitivity-in-private-data-analysis/>
- [21] Google SRE Book (2016). Service Level Objectives and Error Budgets. <https://sre.google/sre-book/service-level-objectives/>
- [22] Enabling Mission-Critical Communication via VoLTE for Public Safety Networks - Varinder Kumar Sharma - IJAIDR Volume 10, Issue 1, January-June 2019. DOI 10.71097/IJAIDR.v10.i1.1539
- [23] Thallam, N. S. T. (2020). Comparative Analysis of Data Warehousing Solutions: AWS Redshift vs. Snowflake vs. Google BigQuery. *European Journal of Advances in Engineering and Technology*, 7(12), 133-141.
- [24] The Role of Zero-Emission Telecom Infrastructure in Sustainable Network Modernization - Varinder Kumar Sharma - IJFMR Volume 2, Issue 5, September-October 2020. <https://doi.org/10.36948/ijfmr.2020.v02i05.54991>
- [25] Krishna Chaitanaya Chittoor, "Architecting Scalable Ai Systems For Predictive Patient Risk", INTERNATIONAL JOURNAL OF CURRENT SCIENCE, 11(2), PP-86-94, 2021, <https://rjpn.org/ijcspub/papers/IJCSP21B1012.pdf>
- [26] Thallam, N. S. T. (2021). Performance Optimization in Big Data Pipelines: Tuning EMR, Redshift, and Glue for Maximum Efficiency.